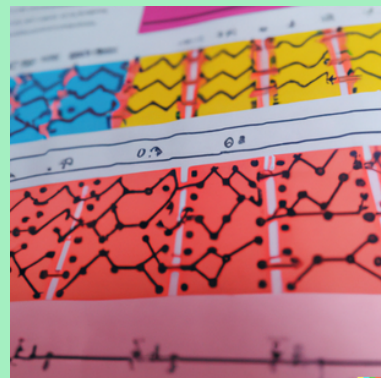# *CASIMIR : un Corpus d'Articles Scientifiques Intégrant les ModIfications et Révisions des auteurs*

Léane Jourdan, Florian Boudin, Richard Dufour, Nicolas Hernandez

{prénom.nom}@univ-nantes.fr

# CONTEXTE

## *Motivations*

- Écrire un article scientifique est une tâche difficile
- De bonnes compétences rédactionnelles sont indispensables
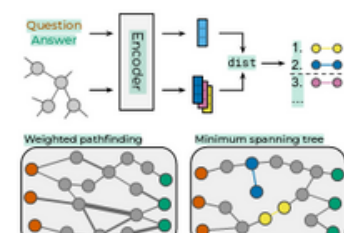- Particulièrement pour les jeunes chercheurs et les non natifs anglophones

## *Objectif*

Offrir une ressource sur la révision d'articles scientifiques



Exemple de révisions

# CASIMIR : un Corpus d'Articles Scientifiques Intégrant les ModIfications et Révisions des auteurs
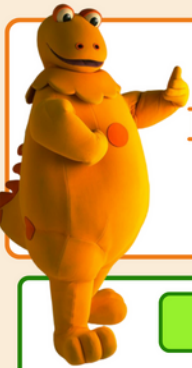
Léane Jourdan, Florian Boudin, Richard Dufour, Nicolas Hernandez

{prénom.nom}@univ-nantes.fr

## Contenu de l'article

- Motivations
- Contenu du corpus
- Processus de collecte

## Prochaines étapes et difficultés

- Conversion de PDF
- Alignement de documents
- Extraction de révisions
- Annotation des types de révisions